

© 2000 Oxford University Press

SURVEY AND SUMMARY

ADEPTs: information necessary for subcellular distribution of eukaryotic sorting isozymes resides in domains missing from eubacterial and archaeal counterparts

David R. Stanford, Nancy C. Martin¹ and Anita K. Hopper*

Department of Biochemistry and Molecular Biology, Pennsylvania State University College of Medicine, H171, 500 University Drive, Hershey, PA 17033, USA and ¹Department of Biochemistry, University of Louisville School of Medicine, 312 Abraham Flexner Way, Louisville, KY 40202, USA

Received September 7, 1999; Revised and Accepted November 22, 1999

ABSTRACT

Sorting isozymes are encoded by single genes, but the encoded proteins are distributed to multiple subcellular compartments. We surveyed the predicted protein sequences of several nucleic acid interacting sorting isozymes from the eukaryotic taxonomic domain and compared them with their homologs in the archaeal and eubacterial domains. Here, we summarize the data showing that the eukaryotic sorting isozymes often possess sequences not present in the archaeal and eubacterial counterparts and that the additional sequences can act to target the eukaryotic proteins to their appropriate subcellular locations. Therefore, we have named these protein domains ADEPTs (Additional Domains for Eukaryotic Protein Targeting). Identification of additional domains by phylogenetic comparisons should be generally useful for locating candidate sequences important for subcellular distribution of eukaryotic proteins.

INTRODUCTION

Eukaryotes are typified by the possession of organelles, generating numerous subcellular locations separated from one another by one or more membranes. Generally the different subcellular compartments carry out unique biochemical reactions. However, sometimes the same catalytic activity is found in more than one subcellular compartment. There are three different mechanisms used by eukaryotic cells to deliver the same enzymatic activity to more than one subcellular location. First, the same catalytic activity may be encoded by dissimilar genes. For example, cognate mitochondrial and cytosolic aminoacyl-tRNA synthetases can be quite distinct (1,2). Second, a catalytic activity may be encoded by multiple similar genes, each encoding an isozyme with unique subcellular distribution.

The yeast genes, *ADH1*, *ADH2* and *ADH3*, provide an example of this type of mechanism (3). Finally, a single gene may encode two or more isozymes with different subcellular distributions. These proteins are called 'sorting isozymes' and are involved in many important metabolic processes (for a review see 4,5).

Sorting isozymes must contain information necessary for protein distribution to different compartments without compromising catalytic activity. Cellular mechanisms that achieve this are varied. In some cases, alternative transcriptional initiation generates mRNAs that encode the catalytic portion with or without signals for specific compartments. In other cases, the same end is achieved by alternative translational initiation or alternative splicing. Finally post-translational modifications can also alter the targeting information without altering catalytic activity (for a review see 4,5). In this report we focus on the cis-acting signals responsible for sorting isozyme distribution.

Genome sequencing efforts have generated information for several archaeal (six are complete and a few others are nearing completion: TIGR, <http://www.tigr.org/tdb/mdb/mdh.html>), many eubacterial (19 are complete and many others are well underway), many, many viral and several eukaryotic nuclear as well as over 100 mitochondrial and 11 chloroplast organellar genomes (see Entrez Genomes at NCBI, <http://www.ncbi.nlm.nih.gov/Entrez/Genome/org.html>). Indeed, the sequences of two eukaryotic nuclear genomes are virtually complete (6,7). If one assumes that sequences important to catalytic function will be conserved, then comparisons of eukaryotic sorting isozymes to their counterpart proteins in non-eukaryotic organisms might reveal the regions of the proteins serving the sorting function.

To test this assumption we conducted phylogenetic comparisons of five proteins. We chose genes that had been functionally characterized by cell biology and molecular biology experiments for their nuclear and mitochondrial targeting signals and some for cytoplasmic retention/nuclear export signals. We used three criteria to choose those proteins. (i) At least one eukaryotic member of the family has been shown directly to be a sorting isozyme and there is detailed information regarding the cis-acting sequences involved in subcellular distribution

*To whom correspondence should be addressed. Tel: +1 717 531 6008; Fax: +1 717 531 7072; Email: ahopper@psu.edu

Table 1. Accession numbers

Organism	Mod5p/Mod5A	Tnn1p	His1p/His1A	Ccp1p	Ung1p	Accession
<i>Escherichia coli</i>	AF010371	AF010372	AF010373	AF010374	AF010375	AF010376
<i>Salmonella typhimurium</i>	AF010377	AF010378	AF010379	AF010380	AF010381	AF010382
<i>Yersinia enterocolitica</i>	AF010383	AF010384	AF010385	AF010386	AF010387	AF010388
<i>Staphylococcus aureus</i>	AF010389	AF010390	AF010391	AF010392	AF010393	AF010394
<i>Streptococcus pneumoniae</i>	AF010395	AF010396	AF010397	AF010398	AF010399	AF010400
<i>Neisseria meningitidis</i>	AF010401	AF010402	AF010403	AF010404	AF010405	AF010406
<i>Haemophilus influenzae</i>	AF010407	AF010408	AF010409	AF010410	AF010411	AF010412
<i>Legionella pneumophila</i>	AF010413	AF010414	AF010415	AF010416	AF010417	AF010418
<i>Campylobacter jejuni</i>	AF010419	AF010420	AF010421	AF010422	AF010423	AF010424
<i>Shigella flexneri</i>	AF010425	AF010426	AF010427	AF010428	AF010429	AF010430
<i>Escherichia coli</i>	AF010431	AF010432	AF010433	AF010434	AF010435	AF010436
<i>Salmonella typhimurium</i>	AF010437	AF010438	AF010439	AF010440	AF010441	AF010442
<i>Yersinia enterocolitica</i>	AF010443	AF010444	AF010445	AF010446	AF010447	AF010448
<i>Staphylococcus aureus</i>	AF010449	AF010450	AF010451	AF010452	AF010453	AF010454
<i>Streptococcus pneumoniae</i>	AF010455	AF010456	AF010457	AF010458	AF010459	AF010460
<i>Neisseria meningitidis</i>	AF010461	AF010462	AF010463	AF010464	AF010465	AF010466
<i>Haemophilus influenzae</i>	AF010467	AF010468	AF010469	AF010470	AF010471	AF010472
<i>Legionella pneumophila</i>	AF010473	AF010474	AF010475	AF010476	AF010477	AF010478
<i>Campylobacter jejuni</i>	AF010479	AF010480	AF010481	AF010482	AF010483	AF010484
<i>Shigella flexneri</i>	AF010485	AF010486	AF010487	AF010488	AF010489	AF010490
<i>Escherichia coli</i>	AF010491	AF010492	AF010493	AF010494	AF010495	AF010496
<i>Salmonella typhimurium</i>	AF010497	AF010498	AF010499	AF010500	AF010501	AF010502
<i>Yersinia enterocolitica</i>	AF010503	AF010504	AF010505	AF010506	AF010507	AF010508
<i>Staphylococcus aureus</i>	AF010509	AF010510	AF010511	AF010512	AF010513	AF010514
<i>Streptococcus pneumoniae</i>	AF010515	AF010516	AF010517	AF010518	AF010519	AF010520
<i>Neisseria meningitidis</i>	AF010521	AF010522	AF010523	AF010524	AF010525	AF010526
<i>Haemophilus influenzae</i>	AF010527	AF010528	AF010529	AF010530	AF010531	AF010532
<i>Legionella pneumophila</i>	AF010533	AF010534	AF010535	AF010536	AF010537	AF010538
<i>Campylobacter jejuni</i>	AF010539	AF010540	AF010541	AF010542	AF010543	AF010544
<i>Shigella flexneri</i>	AF010545	AF010546	AF010547	AF010548	AF010549	AF010550
<i>Escherichia coli</i>	AF010551	AF010552	AF010553	AF010554	AF010555	AF010556
<i>Salmonella typhimurium</i>	AF010557	AF010558	AF010559	AF010560	AF010561	AF010562
<i>Yersinia enterocolitica</i>	AF010563	AF010564	AF010565	AF010566	AF010567	AF010568
<i>Staphylococcus aureus</i>	AF010569	AF010570	AF010571	AF010572	AF010573	AF010574
<i>Streptococcus pneumoniae</i>	AF010575	AF010576	AF010577	AF010578	AF010579	AF010580
<i>Neisseria meningitidis</i>	AF010581	AF010582	AF010583	AF010584	AF010585	AF010586
<i>Haemophilus influenzae</i>	AF010587	AF010588	AF010589	AF010590	AF010591	AF010592
<i>Legionella pneumophila</i>	AF010593	AF010594	AF010595	AF010596	AF010597	AF010598
<i>Campylobacter jejuni</i>	AF010599	AF010600	AF010601	AF010602	AF010603	AF010604
<i>Shigella flexneri</i>	AF010605	AF010606	AF010607	AF010608	AF010609	AF010610
<i>Escherichia coli</i>	AF010611	AF010612	AF010613	AF010614	AF010615	AF010616
<i>Salmonella typhimurium</i>	AF010617	AF010618	AF010619	AF010620	AF010621	AF010622
<i>Yersinia enterocolitica</i>	AF010623	AF010624	AF010625	AF010626	AF010627	AF010628
<i>Staphylococcus aureus</i>	AF010629	AF010630	AF010631	AF010632	AF010633	AF010634
<i>Streptococcus pneumoniae</i>	AF010635	AF010636	AF010637	AF010638	AF010639	AF010640
<i>Neisseria meningitidis</i>	AF010641	AF010642	AF010643	AF010644	AF010645	AF010646
<i>Haemophilus influenzae</i>	AF010647	AF010648	AF010649	AF010650	AF010651	AF010652
<i>Legionella pneumophila</i>	AF010653	AF010654	AF010655	AF010656	AF010657	AF010658
<i>Campylobacter jejuni</i>	AF010659	AF010660	AF010661	AF010662	AF010663	AF010664
<i>Shigella flexneri</i>	AF010665	AF010666	AF010667	AF010668	AF010669	AF010670
<i>Escherichia coli</i>	AF010671	AF010672	AF010673	AF010674	AF010675	AF010676
<i>Salmonella typhimurium</i>	AF010677	AF010678	AF010679	AF010680	AF010681	AF010682
<i>Yersinia enterocolitica</i>	AF010683	AF010684	AF010685	AF010686	AF010687	AF010688
<i>Staphylococcus aureus</i>	AF010689	AF010690	AF010691	AF010692	AF010693	AF010694
<i>Streptococcus pneumoniae</i>	AF010695	AF010696	AF010697	AF010698	AF010699	AF010700
<i>Neisseria meningitidis</i>	AF010701	AF010702	AF010703	AF010704	AF010705	AF010706
<i>Haemophilus influenzae</i>	AF010707	AF010708	AF010709	AF010710	AF010711	AF010712
<i>Legionella pneumophila</i>	AF010713	AF010714	AF010715	AF010716	AF010717	AF010718
<i>Campylobacter jejuni</i>	AF010719	AF010720	AF010721	AF010722	AF010723	AF010724
<i>Shigella flexneri</i>	AF010725	AF010726	AF010727	AF010728	AF010729	AF010730
<i>Escherichia coli</i>	AF010731	AF010732	AF010733	AF010734	AF010735	AF010736
<i>Salmonella typhimurium</i>	AF010737	AF010738	AF010739	AF010740	AF010741	AF010742
<i>Yersinia enterocolitica</i>	AF010743	AF010744	AF010745	AF010746	AF010747	AF010748
<i>Staphylococcus aureus</i>	AF010749	AF010750	AF010751	AF010752	AF010753	AF010754
<i>Streptococcus pneumoniae</i>	AF010755	AF010756	AF010757	AF010758	AF010759	AF010760
<i>Neisseria meningitidis</i>	AF010761	AF010762	AF010763	AF010764	AF010765	AF010766
<i>Haemophilus influenzae</i>	AF010767	AF010768	AF010769	AF010770	AF010771	AF010772
<i>Legionella pneumophila</i>	AF010773	AF010774	AF010775	AF010776	AF010777	AF010778
<i>Campylobacter jejuni</i>	AF010779	AF010780	AF010781	AF010782	AF010783	AF010784
<i>Shigella flexneri</i>	AF010785	AF010786	AF010787	AF010788	AF010789	AF010790
<i>Escherichia coli</i>	AF010791	AF010792	AF010793	AF010794	AF010795	AF010796
<i>Salmonella typhimurium</i>	AF010797	AF010798	AF010799	AF010800	AF010801	AF010802
<i>Yersinia enterocolitica</i>	AF010803	AF010804	AF010805	AF010806	AF010807	AF010808
<i>Staphylococcus aureus</i>	AF010809	AF010810	AF010811	AF010812	AF010813	AF010814
<i>Streptococcus pneumoniae</i>	AF010815	AF010816	AF010817	AF010818	AF010819	AF010820
<i>Neisseria meningitidis</i>	AF010821	AF010822	AF010823	AF010824	AF010825	AF010826
<i>Haemophilus influenzae</i>	AF010827	AF010828	AF010829	AF010830	AF010831	AF010832
<i>Legionella pneumophila</i>	AF010833	AF010834	AF010835	AF010836	AF010837	AF010838
<i>Campylobacter jejuni</i>	AF010839	AF010840	AF010841	AF010842	AF010843	AF010844
<i>Shigella flexneri</i>	AF010845	AF010846	AF010847	AF010848	AF010849	AF010850
<i>Escherichia coli</i>	AF010851	AF010852	AF010853	AF010854	AF010855	AF010856
<i>Salmonella typhimurium</i>	AF010857	AF010858	AF010859	AF010860	AF010861	AF010862
<i>Yersinia enterocolitica</i>	AF010863	AF010864	AF010865	AF010866	AF010867	AF010868
<i>Staphylococcus aureus</i>	AF010869	AF010870	AF010871	AF010872	AF010873	AF010874
<i>Streptococcus pneumoniae</i>	AF010875	AF010876	AF010877	AF010878	AF010879	AF010880
<i>Neisseria meningitidis</i>	AF010881	AF010882	AF010883	AF010884	AF010885	AF010886
<i>Haemophilus influenzae</i>	AF010887	AF010888	AF010889	AF010890	AF010891	AF010892
<i>Legionella pneumophila</i>	AF010893	AF010894	AF010895	AF010896	AF010897	AF010898
<i>Campylobacter jejuni</i>	AF010899	AF010900	AF010901	AF010902	AF010903	AF010904
<i>Shigella flexneri</i>	AF010905	AF010906	AF010907	AF010908	AF010909	AF010910
<i>Escherichia coli</i>	AF010911	AF010912	AF010913	AF010914	AF010915	AF010916
<i>Salmonella typhimurium</i>	AF010917	AF010918	AF010919	AF010920	AF010921	AF010922
<i>Yersinia enterocolitica</i>	AF010923	AF010924	AF010925	AF010926	AF010927	AF010928
<i>Staphylococcus aureus</i>	AF010929	AF010930	AF010931	AF010932	AF010933	AF010934
<i>Streptococcus pneumoniae</i>	AF010935	AF010936	AF010937	AF010938	AF010939	AF010940
<i>Neisseria meningitidis</i>	AF010941	AF010942	AF010943	AF010944	AF010945	AF010946
<i>Haemophilus influenzae</i>	AF010947	AF010948	AF010949	AF010950	AF010951	AF010952
<i>Legionella pneumophila</i>	AF010953	AF010954	AF010955	AF010956	AF010957	AF010958
<i>Campylobacter jejuni</i>	AF010959	AF010960	AF010961	AF010962	AF010963	AF010964
<i>Shigella flexneri</i>	AF010965	AF010966	AF010967	AF010968	AF010969	AF010970
<i>Escherichia coli</i>	AF010971	AF010972	AF010973	AF010974	AF010975	AF010976
<i>Salmonella typhimurium</i>	AF010977	AF010978	AF010979	AF010980	AF010981	AF010982
<i>Yersinia enterocolitica</i>	AF010983	AF010984	AF010985	AF010986	AF010987	AF010988
<i>Staphylococcus aureus</i>	AF010989	AF010990	AF010991	AF010992	AF010993	AF010994
<i>Streptococcus pneumoniae</i>	AF010995	AF010996	AF010997	AF010998	AF010999	AF011000

Organism	Mod5p/Mod5A	Tnn1p	His1p/His1A	Ccp1p	Ung1p	Accession
<i>Escherichia coli</i>	AF010371	AF010372	AF010373	AF010374	AF010375	AF010376
<i>Salmonella typhimurium</i>	AF010377	AF010378	AF010379	AF010380	AF010381	AF010382
<i>Yersinia enterocolitica</i>	AF010383	AF010384	AF010385	AF010386	AF010387	AF010388
<i>Staphylococcus aureus</i>	AF010389	AF010390	AF010391	AF010392	AF010393	AF010394
<i>Streptococcus pneumoniae</i>	AF010395	AF010396	AF010397	AF010398	AF010399	AF010400
<i>Neisseria meningitidis</i>	AF010401	AF010402	AF010403	AF010404	AF010405	AF010406
<i>Haemophilus influenzae</i>	AF010407	AF010408	AF010409	AF010410	AF010411	AF010412
<i>Legionella pneumophila</i>	AF010413	AF010414	AF010415	AF010416	AF010417	AF010418
<i>Campylobacter jejuni</i>	AF010419	AF010420	AF010421	AF010422	AF010423	AF010424
<i>Shigella flexneri</i>	AF010425	AF010426	AF010427	AF010428	AF010429	AF010430
<i>Escherichia coli</i>	AF010431	AF010432	AF010433	AF010434	AF010435	AF010436
<i>Salmonella typhimurium</i>	AF010437	AF010438	AF010439	AF010440	AF010441	AF010442
<i>Yersinia enterocolitica</i>	AF010443	AF010444	AF010445	AF010446	AF010447	AF010448
<i>Staphylococcus aureus</i>	AF010449	AF010450	AF010451	AF010452	AF010453	AF010454
<i>Streptococcus pneumoniae</i>	AF010455	AF010456	AF010457	AF010458	AF010459	AF010460
<i>Neisseria meningitidis</i>	AF010461	AF010462	AF010463	AF010464	AF010465	AF010466
<i>Haemophilus influenzae</i>	AF010467	AF010468	AF010469	AF010470	AF010471	AF010472
<i>Legionella pneumophila</i>	AF010473	AF010474	AF010475	AF010476	AF010477	AF010478
<i>Campylobacter jejuni</i>	AF010479	AF010480	AF010481	AF010482	AF010483	AF010484
<i>Shigella flexneri</i>	AF010485	AF010486	AF010487	AF010488	AF010489	AF010490
<i>Escherichia coli</i>	AF010491	AF010492	AF010493	AF010494	AF010495	AF010496
<i>Salmonella typhimurium</i>	AF010497	AF010498	AF010499	AF010500	AF010501	AF010502
<i>Yersinia enterocolitica</i>	AF010503	AF010504	AF010505	AF010506	AF010507	AF010508
<i>Staphylococcus aureus</i>	AF010509	AF010510	AF010511	AF010512	AF010513	AF010514

[illegible]

The data are presented in two ways. Figures S1–S5 available as Supplementary Material at NAR Online, show the actual amino acid sequence alignment information. A score of ≥1 from the BLOSUM 62 matrix is designated as similar while a score of 0 is considered a weak similarity. Amino acids are grouped and colored as follows: aromatic amino acids phenylalanine, tyrosine and tryptophan (FYW) are magenta; hydrophobic amino acids isoleucine, valine, leucine and methionine (IVLM) are cyan; charged/polar amino acids aspartic acid, glutamic acid, glutamine, lysine, arginine, asparagine and histidine (DEQKRNH) are red; small amino acids glycine,

[illegible]

Figures 2-6 show schematic diagrams of the protein alignments based on the sequence alignments described above. Blocks of similar color represent blocks of sequence similarity and are not a representation of any structural information. Different colored boxes represent uninterrupted regions of

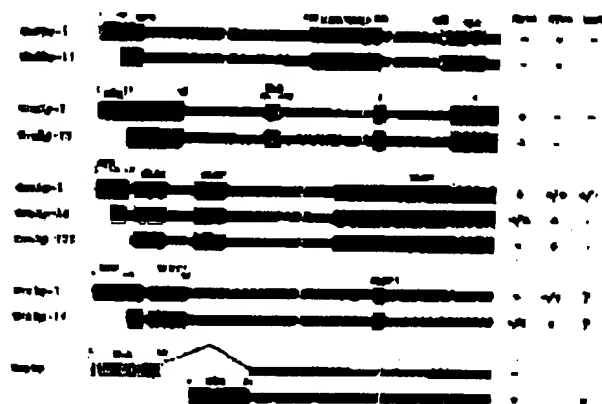


Figure 1. Location of information for subcellular distribution of sorting isozymes. Known and presumed targeting signals are represented as colored boxes. Magenta boxes represent known mitochondrial targeting information. Teal boxes and blue boxes represent known and presumed (NLS?) nuclear targeting information, respectively. Purple boxes may target Trnlp to a subnuclear location and the green boxes in ModSp may be responsible for the predominantly cytosolic distribution of this protein. CRD, cytoplasmic retention domain; NES, nuclear export signal. The black lines represent the conserved regions of each protein and are not to scale. The subcellular distributions of the various forms of each protein are also indicated. For Htr1p, -/+ refers to locations detected upon protein over-expression.

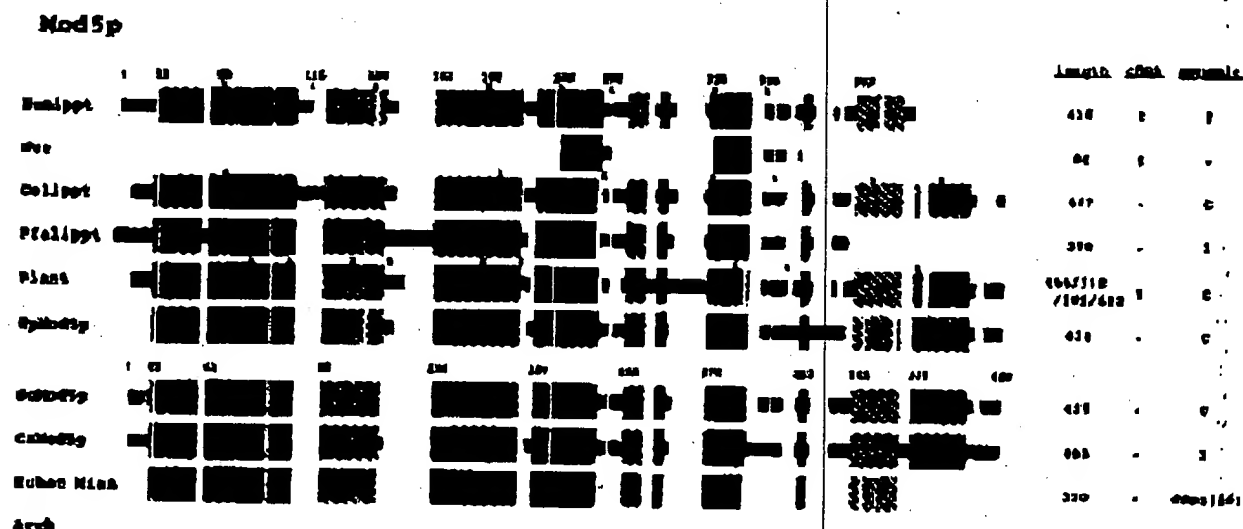
similarity (at least 35%) between the proteins from different organisms. Black lines represent eukaryotic sequences not generally similar to each other. Gray lines represent prokaryotic sequences not generally similar to each other or the eukaryotic sequences. Not all the sequences depicted are complete and

some of the eukaryotic peptides judged to be too incomplete are not shown in the schematic diagrams. Eight eukaryotes were selected to represent the domain Eukarya: *Homo sapiens*, *Mus musculus*, *Caenorhabditis elegans*, *Plasmodium falciparum*, *Schizosaccharomyces pombe*, *Saccharomyces cerevisiae* and *Candida albicans*. Plants are usually represented as a composite diagram due to the lack of complete sequence information. An I to the right of the schematics designates incomplete information and n C designates complete cDNA or genomic DNA sequence information. The lengths of the polypeptide chains are indicated and where a composite schematic is shown the lengths of the individual polypeptide chains are separated by slashes. The eubacterial and archaeal schematics are derived from consensus sequences and the number of peptides used to generate the consensus is also indicated. Where information is available concerning the site of intron-exon junctions, the locations of introns are marked with an x.

RESULTS AND DISCUSSION

ModSp homologs and conservation of regions for subcellular distribution

We previously reported an alignment of ModSp/MiaA from 33 eubacteria and three eukaryotes (13). Our continued search for ModSp homologs has now uncovered ModSp/MiaA in 45 eubacteria (see Table 1). Two eubacterial organisms do not contain a miaA gene (*Mycoplasma genitalium* and *Mycoplasma pneumoniae*) while one, *Porphyromonas gingivalis*, contains two miaA genes. Seventeen eukaryotic homologs were identified in fifteen organisms (*H.sapiens*, *M.murinus*, *Drosophila melanogaster*, *C.elegans*, *P.falciparum*, *Cryptosporidium parvum*, *Leishmania major*, *Trypanosoma brucei*, *Arabidopsis thaliana*, *Oryza sativa*, *S.pombe*, *S.cerevisiae*, *C.albicans*,



Kluyveromyces fragilis and *Neurospora crassa*). *Saccharomyces cerevisiae* and *C. elegans* have only one gene encoding this protein. Only eight of the eukaryotic Mod5s are shown in Figure 2 and the 46 eubacterial MiaA homologs are represented in Figure 2 as a consensus schematic. The entry for plants in Figure 2 represents a composite of three *A. thaliana* homologs and one homolog from rice. No homologs were identified in archaea, consistent with the fact that i⁶A has not been found on tRNAs isolated from organisms in the archaeal domain (14,15).

By alternative translational starts the *S. cerevisiae* MOD5 gene encodes two proteins, Mod5p-I and Mod5p-II (16), which are differentially partitioned between the cytoplasm, mitochondria and nucleus (17). Mod5p-I is located in the mitochondrial and cytosolic compartments whereas Mod5p-II is in the cytosol and the nucleus. Amino acids 1-20 comprise a mitochondrial targeting sequence (MTS) necessary for distribution of Mod5p-I to the mitochondria (17).

MTSs are usually located at the N-terminus, contain basic and hydrophobic amino acids and are predicted to form amphiphilic α -helices; however, there is no linear consensus sequence for mitochondrial targeting information (18,19). To assess whether other eukaryotes may utilize the same strategy as that for *S. cerevisiae*, we investigated the N-terminal regions of the other eukaryotic Mod5 proteins. Five of the eukaryotic homologs (*S. cerevisiae*, *C. elegans*, *C. albicans*, *P. falciparum* and one of the homologs from *A. thaliana*) contain multiple ATGs at the beginning of the coding region (Fig. 2), while for most of the other eukaryotes there is insufficient information available to predict whether or not multiple translation initiations give rise to different isozymes. The amphiphilic nature of these N-terminal peptides was investigated by plotting them on a helical wheel projection (not shown). In addition to *S. cerevisiae*, the *C. elegans* and *C. albicans* N-terminal regions resemble other MTSs (18,19). Thus, we predict that the *C. elegans* and *C. albicans* Mod5 proteins will also be sorted between the cytoplasm and mitochondria. The N-terminal regions of the *P. falciparum* homolog and the *A. thaliana* homolog with an N-terminal extension (Fig. S1, Atha1p1p) do not resemble other MTSs. In general, the eubacterial proteins do not have this N-terminal extension bolstering the idea that this extra domain found in the eukaryotic proteins is used for targeting.

Arabidopsis thaliana has at least three genes predicted to encode Mod5 proteins; therefore, different genes may well provide the same catalytic activity to different compartments for this organism. While additional information concerning *A. thaliana* and other eukaryotic organisms will be required to determine how mitochondrial/chloroplast/cytoplasmic/nuclear sorting may be achieved, it appears that for the Mod5 family sometimes one gene codes a catalytic activity found in multiple compartments whereas in other cases, two or more genes may code the isozymes.

Nearly all of the eukaryotic Mod5 proteins possess ~50 amino acids at the C-terminus that are not present in the eubacterial MiaA proteins (Fig. 2). The *S. cerevisiae* Mod5p nuclear localization sequence (NLS) maps within this 'additional domain' (amino acids 408-428; 13). In all of the other eukaryotes where sufficient sequence information is available (Fig. 2; *S. pombe*, *C. albicans*, *C. elegans*, rice and one of the *A. thaliana* homologs), the C-terminal region is similar leading to the prediction that they all contain a NLS and that a portion of the

Mod5p pool in these organisms will also be located in the nucleus. Only one of the three *A. thaliana* homologs contains this NLS region while the others lack it (Fig. S1, not shown in Fig. 2), again suggesting that multiple genes encode differently located Mod5p in *A. thaliana*.

Besides the N-terminal and C-terminal additional domains, the eukaryotic Mod5 proteins also contain internal domains not found in the eubacterial homologs (Fig. 2). These internal additions overlap the region between amino acids 240 and 280 that were previously mapped to function in maintenance of the yeast Mod5p cytosolic pool (13). As all the eukaryotic sequences contain a similar region, we predict each of the eukaryotic counterparts also has a cytosolic pool of this protein.

A portion of the *S. cerevisiae* Mod5p-II resides in the nucleolus (13). The information used for nucleolar location has not been mapped. If, like the NLS and MTS, the nucleolar targeting/retention information resides in motifs absent from the eubacterial counterparts, then candidate locations for nucleolar targeting are between amino acids 303 and 345 and/or 373 and 408.

Trm1p homologs and conservation of regions for subcellular distribution

TRM1 genes are found in eukaryotes and archaea, but are generally not present in eubacteria (Fig. 3). In addition to the Trm1p homologs that have already been reported (20,21; six from the archaeal domain, *Aquifex aeolicus*, *S. cerevisiae*, *S. pombe*, *C. elegans* and human) our searches revealed three additional archaeal homologs and incomplete sequences for mouse, rat, zebrafish, *D. melanogaster*, *P. falciparum*, *C. parvum*, *T. brucei*, *A. thaliana*, rice, *Brassica*, *Zea mays* and *C. albicans*. There is only a single eubacterial organism, *A. aeolicus*, that contains a *trm1* gene and this is likely a result of horizontal transfer (22-24). In agreement with our alignments, previous studies of tRNA modification have failed to uncover m²G in eubacterial tRNAs (14,15,25).

Eukaryotic and archaeal Trm1 proteins have considerable sequence similarity. However, like Mod5p, the eukaryotic proteins contain extra sequence information at the N- and C-termini and internally. The *S. cerevisiae* TRM1 gene contains ATG codons at positions 1 and 17. Human Trm1p contains two ATGs within the first 37 codons while mouse Trm1p contains three ATGs within the first 32 codons. Of the eukaryotic genes that have been sequenced at the N-terminus, only two, from *C. elegans* and *D. melanogaster* do not have multiple ATGs within the first 50 codons.

Some mitochondrial tRNAs of *S. cerevisiae* are modified by Trm1p and amino acids 1-48 of the *S. cerevisiae* Trm1p are sufficient to target this protein to mitochondria whereas amino acids 1-16 are not sufficient (26). There are several reports of m²G in mitochondrial and chloroplast tRNAs (27), but unfortunately the TRM1 genes have not been sequenced for the organisms demonstrated to contain m²G modified mitochondrial or chloroplast tRNAs. The N-terminus of the human Trm1p contains no acidic amino acids (Fig. S2) and when projected upon a helical wheel, it is predicted to have an amphiphilic structure, characteristic of MTSs (19). Thus, the human gene could encode a Trm1p that sorts to the mitochondria. The rodent homologs are very similar to the human in this region and the *C. albicans* Trm1p N-terminus contains what appears to be a very good MTS. As the *C. elegans* genome contains only a

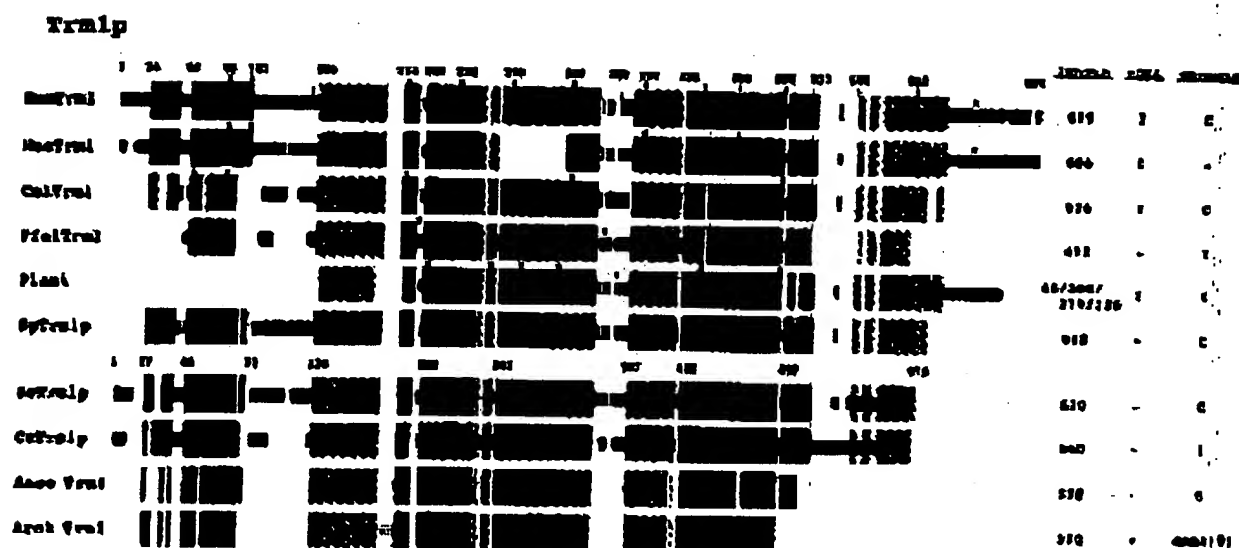


Figure 3. Schematic diagram of Trm1p alignment. A sequence alignment of all identified Trm1p homologs can be found in Figure S2. Nine actual Trm1 peptides were identified and are represented as a consensus sequence. One trm1 homolog was identified in the eubacterial domain. The schematic for plant in this figure is a composite of *A.thaliana*, *O.sativa* and *Brassica*. Regions of uninterrupted sequence similarity are shown as crosshatched colored boxes. See Methods for additional explanations.

single *TRM1* gene, it is likely that this gene provides the mitochondrial pool of tRNA (guanine-26, N^2 - N^2) methyltransferase, if this modification occurs in *C.elegans* mitochondria.

Saccharomyces cerevisiae Trm1p is also targeted to the nucleus and an efficient NLS resides between amino acids 95 and 102 (28). All the other eukaryotic Trm1 proteins contain extra sequence information in this same region (Fig. 3, black region between 103 and 156 of human Trm1p). The *C.elegans* (21) and *D.melanogaster* proteins contain basic amino acids resembling the simple basic type of NLS in this region (see the review in 29), perhaps indicating a functional role in nucleus location. The corresponding extra sequences in human, mouse and *S.pombe* are not nearly as basic as the *S.cerevisiae* Trm1p sequence and neither a simple nor bipartite basic NLS motif can be identified in this region. However, it has recently become apparent that there are multiple nuclear import receptors in eukaryotic cells that have substrate specificities not yet delineated (see the review in 30). If the ADEPT regions of human, mouse and *S.pombe* Trm1p are used to sort this protein to the nucleus, as is the case in *S.cerevisiae*, then phylogenetic comparisons and sequence alignments may be a useful means to delineate non-conventional NLS motifs.

The eukaryotic genes also predict a large C-terminal region and a smaller region (between amino acids 346 and 367 in *S.cerevisiae*) not found in the archaeal proteins (Fig. 3). A zinc finger is present in the eukaryotic proteins (amino acids 348-387 human Trm1p) that is present in only half of the prokaryotic proteins. When present in prokaryotic proteins, the finger loop is much smaller than that found in eukaryotic proteins. The nuclear pool of Trm1p in *S.cerevisiae* is located at the inner surface of the nuclear membrane (28,31). If location at this subnuclear site is achieved via an ADEPT, then we predict that the targeting information will map to either the large C-terminal or

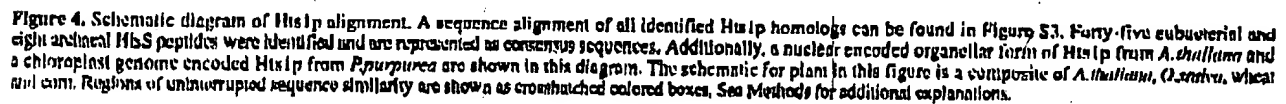
the smaller upstream eukaryotic additional sequences (Fig. 1, purple boxes and Fig. 3).

Others (32) have reported results both consistent and inconsistent with our hypothesis. Deletion of the first 44 amino acids of *S.cerevisiae* Trm1p does not influence enzymatic activity, which is in accord with previous work demonstrating that this region contains targeting information (26) as well as our prediction that this region of the other eukaryotic proteins will supply targeting information. However, a deletion of just five amino acids at the C-terminus of *S.cerevisiae* Trm1p causes a significant reduction in activity (32). This result is inconsistent with our model in that all of the prokaryotic trm1 proteins lack this region and thus it is not expected to influence enzymatic activity. It is conceivable that an alteration in this region of the eukaryotic proteins may effect the higher order structure of the protein and interfere with activity.

His1p homologs and conservation of regions for subcellular distribution

HTS1 encodes histidine-tRNA synthetase, which is known as HisS in prokaryotes. Forty-five eubacterial and eight archaeal homologs were identified and 30 eukaryotic homologs were found. This enzyme is very similar in all three taxonomic domains (Fig. 4). Signature sequences can be identified that distinguish the eubacterial and archaeal proteins, and in some regions the archaeal signature is more similar to that of eukaryotes than to that of eubacteria.

Six of the eukaryotic homologs contain multiple ATGs in their 5' regions. However, the majority of the eukaryotic sequences are incomplete in this region and therefore we are unable to predict whether they encode proteins that differ at the N-terminus. In humans there are two genes arranged head-to-head that code for histidine-tRNA synthetases (Fig. 4,



Both *Xenopus* oocytes (35) and *S.cerevisiae* (36) aminoacylate tRNAs inside the nucleus as well as in the cytosol. Therefore, there must be nuclear pools of aminoacyl-tRNA synthetases. If H1p indeed possesses information that directs it to the nuclear interior, the targeting information could be located in the N-terminal region (Fig. S3, amino acids 20–53 of HumIARS). The additional sequences at this location in eukaryotic proteins contain basic residues resembling conventional NLS motifs (37). Fine mapping of the MTS in this

Organisms in all three domains contain ATP (CTP): tRNA nucleotidyltransferase activity. However, the archaeal Cca proteins differ extensively from the eubacterial and eukaryotic Cca proteins (43). Nevertheless, all possess 'nucleotidyltransferase' motifs. Of the proteins we studied Cca1p is the least well conserved between eubacteria and eukaryotes. Large regions of sequence similarity, as found for the other proteins in our analysis, are lacking in this family. Sixteen eukaryotic

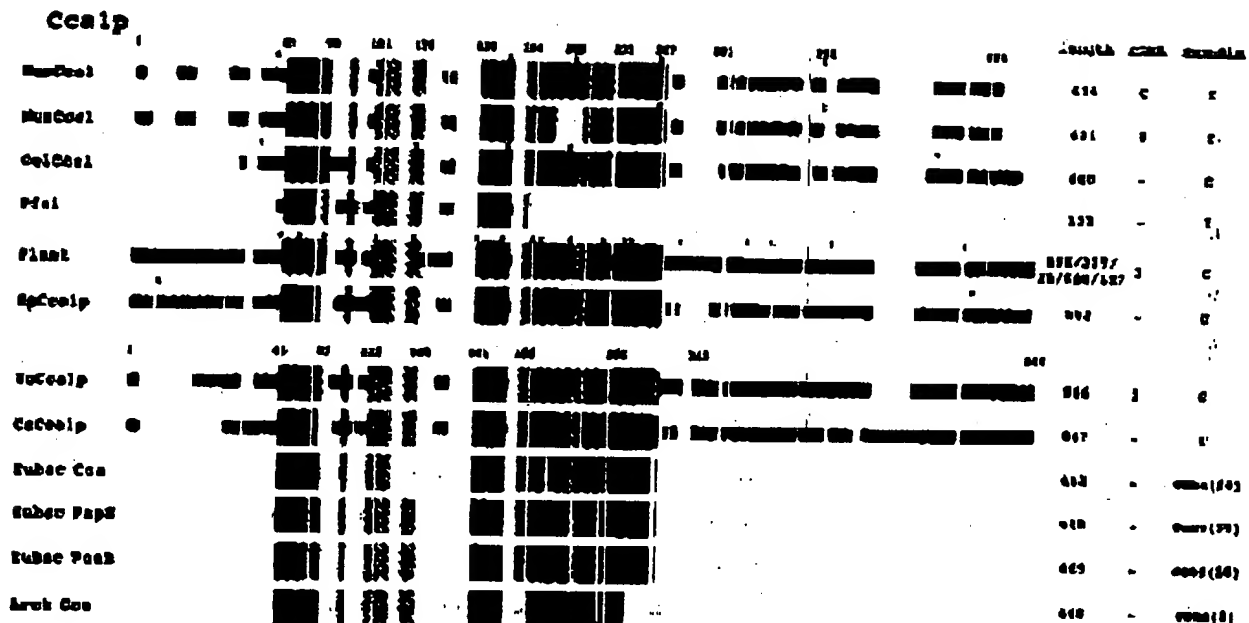


Figure 5. Schematic diagram of Cca1p alignment. A sequence alignment of all identified Cca1p homologs can be found in Figure S4. Eight archael Cca1p peptides were identified and are represented as a consensus schematic. Sixty-five homologs were identified in the eubacterial domain. The eubacterial homologs fall into three classes and a consensus schematic is presented for each class: Cca-14, Pap-12 and PcaB-16. The schematic for plant in this figure is a composite of *A. thaliana*, *O. sativa*, lupine and *G. max*. Regions of uninterrupted sequence similarity are shown as crosshatched colored boxes. See Methods for additional explanations.

homologs were identified in the following organisms: *S. cerevisiae*, *S. pombe*, *C. albicans*, human, mouse, rat, *C. elegans*, *D. melanogaster*, *A. thaliana*, lupine, rice, *Glycine max*, *L. major*, *Brugia malayi* and *P. falciparum*. Eight archael homologs and 65 eubacterial homologs were identified. The latter have been grouped into three classes (Cca, Pap and PcaB) based on the sequence alignments as well as previous nomenclature. A consensus schematic is shown for each of these three classes of eubacterial proteins in Figure 5.

In *S. cerevisiae* the *CCA1* gene encodes three proteins (Cca1p-I, Cca1p-II and Cca1p-III) that result from differential translation starts at three in-frame AUGs (44). Eight of the eukaryotic genes have multiple ATGs in this N-terminal region (Fig. S4), suggesting that multiple forms of Cca1p could also be produced by these genes.

Cca1p-I from *S. cerevisiae* is located primarily in mitochondria whereas Cca1p-II and Cca1p-III are located both in the cytosol and the nucleus (45). Like MudSp, Trm1p and His1p the N-terminus of *S. cerevisiae* Cca1p contains mitochondrial targeting information. For each of the other eukaryotes where there is sufficient information, the eukaryotic Cca1p counterparts have an N-terminal extension that is absent or different in the eubacterial and archael proteins. This region most likely directs the non-plant Cca1p to mitochondria. Plant Cca1p should also be directed to the chloroplast. As chloroplast targeting information also is usually located at the N-terminus and resembles mitochondrial targeting information (46; for a review see 47), it is difficult to predict the function of the plant N-terminal Cca1p extensions.

Also, since no plant genome has been completely sequenced there could be different genes for mitochondrial and chloroplast CCA activities.

The location of other targeting information for Cca1p is unknown, but there are other regions that contain additions not found in eubacteria (94–103; 109–114 *S. cerevisiae* numbering). There are also extensive regions of the proteins that are dissimilar between eukaryotes and prokaryotes (Fig. 5) that could contain nuclear targeting information.

Ung1p homologs and conservation of regions for subcellular distribution

Uracil-DNA glycosylase (UNG or UDG) is a DNA repair enzyme. The *ung* gene is found in 33 eubacteria, but is not present in archaea. Thus, either another gene product supplies this function or this function is not required. Interestingly, of the 19 complete eubacterial genomes, the *ung* gene is absent from six (*Rickettsia prowazekii*, *Clostridium acetobutylicum*, *Treponema pallidum*, *A. aeolicus*, *Thermotoga maritima* and *Synechocystis*), again suggesting that this function may not be required. Also of note is that within the genus *Clostridium* one organism, *Clostridium difficile*, contains a *ung* gene while *C. acetobutylicum* does not. *UNG* genes are also present in some viruses and consensus sequences for the Ung protein from 23 Herpes simplex viruses and five pox viruses are shown in Figure S5.

The human homolog of this enzyme is the most thoroughly studied. BLAST searches revealed Ung homologs in 11 other

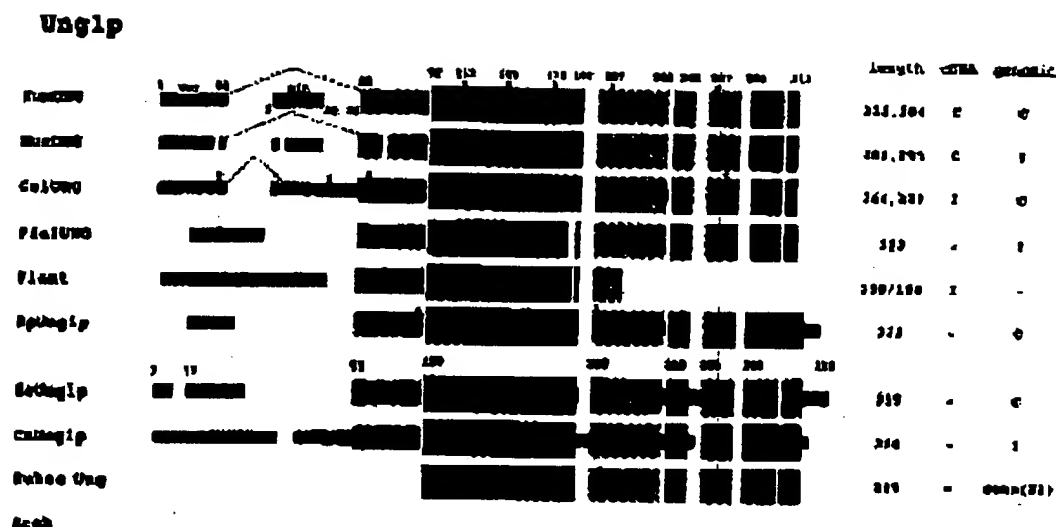


Figure 6. Schematic diagram of Ung1p alignment. A sequence alignment of all identified Ung1p homologs can be found in Figure S5. Ung1p was not identified in the archaeal domain. Thirty-three homologs were identified in the eubacterial domain and a consensus schematic is presented for these homologs. The schematic for plant in this figure is a composite of poplar and yamnia. Regions of uninterrupted sequence similarity are shown as crosshatched colored boxes. Alternatively spliced exons are also indicated. See Methods for additional explanations.

eukaryotes. The mouse homolog is very similar to the human (90% similarity) and both sort this enzyme between the nucleus and mitochondria via a mechanism that depends on alternative splicing (48,49; Fig. 6). This mechanism may also be used in *C. elegans* as there is an extra 'exon' upstream of the *UNG* gene which could be used to supply additional targeting information. However, this putative exon does not resemble known MTS or NLS motifs. Disregarding this putative exon the *C. elegans* ORF contains four in-frame ATGs. Downstream of AUG2 there is a sequence resembling a MTS, but we were unable to identify a classical simple or bipartite-like NLS in the N-terminal region. In *S. cerevisiae* there are four methionines within the first 50 amino acids and alternative transcription or translation start sites could provide the sorting mechanism for this enzyme; however, the available data (50; P.Burgers, personal communication) indicate that Ung1p is solely nuclear and unlikely to sort to mitochondria in yeast.

Since Ung1p should function within the nucleus of eukaryotes, there should be information to target this enzyme to the nucleus. Most of the eukaryotic and viral Ung proteins contain extra N-terminal sequence information not found in the bacterial counterparts. The human and mouse nuclear targeting information resides within this region and *S. cerevisiae* and *P. falciparum* appear to contain conventional bipartite NLSs within this region.

CONCLUSIONS

We surveyed five families of proteins containing at least one confirmed sorting isozyme. Four of these protein families have members that are highly conserved across taxonomic domains and the eukaryotic proteins contain additional sequences not

found in the eubacterial or archaeal counterparts. Although the fifth protein, Ccd1p, fits the pattern established by the other proteins in a limited sense, large portions of this protein are dissimilar when compared across taxonomic domains.

Additional information can be located at the N- or C-termini or it can be located internally. The location of additional sequence information is conserved, but the sequences are not necessarily similar. It has been proposed that intron locations correspond to positions separating independent functional domains of proteins (51,52). Although our data set is limited, our analysis does not appear to support this view. In general, ADEPTs do not correspond to genomic spliced regions.

We summarize the evidence that the additional sequences can encode information to sort the isozymes to appropriate subcellular locations (Fig. 1). The data lead us to propose the ADEPT hypothesis that similarly located extra information in other eukaryotic homologs will serve the same roles in protein subcellular distribution. We present this type of analysis as a predictive tool. Our results suggest that phylogenetic comparison/multiple sequence alignment will be a useful tool for predicting the cell biological information content of protein sequences. Future mechanistic tests of the sequences identified here will be necessary to determine how accurate these predictions are. However, data to date are quite consistent with the ADEPT concept.

SUPPLEMENTARY MATERIAL

See Supplementary Material available at NAR Online. Update to the published Supplementary Material will be available at <http://www.collmed.psu.edu/labs/ahopper/DRS/ADEPTs/sortpaper.htm>

ACKNOWLEDGEMENTS

We would like to express our gratitude to all members of the various genome sequencing projects (TIGR, Sanger Centre, OU-ACGT, Washington University GSC, BSNR, ChGP, Diversa, GTC, PGP, ASTRA, KDRL, NITE, Genome Therapeutics, Stanford University, UC Berkeley, University of Heidelberg, Uppsala University, University of Minnesota CBC, University of Wisconsin and the Utah Genome Center) for making information available and to the various funding agencies involved with these projects. Web links to the respective sequencing center are provided for each entry in Table S1. We would also like to thank all the members of the Hopper laboratory for their thoughtful discussions regarding this work. This work was supported by grants from the National Science Foundation to A.K.H. and N.C.M. and from the National Institutes of Health to A.K.H.

REFERENCES

- Schneller, J.M., Schneller, C. and Stuhl, A.J. (1978) *Biochem. Biophys. Res. Commun.*, **85**, 1392-1399.
- Tzagoulis, A. and Shianko, A. (1995) *Eur. J. Biochem.*, **230**, 583-586.
- Pilgrim, D. and Young, E.T. (1990) *Mol. Cell. Biol.*, **7**, 294-304.
- Martin, N.C. and Hopper, A.K. (1994) *Biochimie*, **76**, 1161-1167.
- Danpure, C.J. (1995) *Trends Cell Biol.*, **5**, 231-237.
- Goffeau, A., Aert, R., Agostini-Carbone, M., Altieri, A., Aigle, M., Algerghini, L., Altmann, K., Albers, M., Alden, M., Alexandrak, D. et al. (1997) *Nature*, **387** (suppl.), 5-105.
- The *C. elegans* Sequencing Consortium (1998) *Science*, **282**, 2012-2018.
- Altshul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J. (1997) *Nucleic Acids Res.*, **25**, 3389-3402.
- Higgins, D.G., Thompson, J.D. and Gibson, T.J. (1996) *Methods Enzymol.*, **266**, 383-402.
- Thompson, J.D., Gibson, T.J., Plewniak, F., Jeanmougin, P. and Higgins, D.G. (1997) *Nucleic Acids Res.*, **25**, 4876-4882.
- Henikoff, S. and Henikoff, J.G. (1992) *Proc. Natl Acad. Sci. USA*, **89**, 10915-10919.
- Livingstone, C.D. and Barton, G.J. (1996) *Methods Enzymol.*, **266**, 497-512.
- Tolcico, L.H., Bonko, A.I., Arls, J.P., Stanford, D.R., Martin, N.C. and Hopper, A.K. (1999) *Genetics*, **151**, 57-75.
- Maden, H.F.H. (1998) In Grosjean, H. and Benne, R. (eds), *Modification and Editing of RNA*. ASM Press, Washington, DC, pp. 421-440.
- Winkler, M.E. (1998) In Grosjean, H. and Benne, R. (eds), *Modification and Editing of RNA*. ASM Press, Washington, DC, pp. 441-468.
- Gillman, E.C., Slusher, I.B., Martin, N.C. and Hopper, A.K. (1991) *Mol. Cell. Biol.*, **11**, 2382-2390.
- Boguta, M., Hopter, L.A., Shen, W.-C., Gillman, E.C., Martin, N.C. and Hopper, A.K. (1994) *Mol. Cell. Biol.*, **14**, 2298-2306.
- Anzard, O. and Schatz, G. (1988) *Annu. Rev. Cell. Biol.*, **4**, 289-333.
- von Heijne, G. (1986) *EMBO J.*, **5**, 1335-1342.
- Constantinico, F., Bonaschewitz, N., Morfin, Y. and Grosjean, H. (1998) *Nucleic Acids Res.*, **26**, 3753-3761.
- Liu, J., Zhou, Q. and Stray, K.B. (1999) *Gene*, **226**, 73-81.
- Syvanen, M. (1994) *Annu. Rev. Genet.*, **28**, 237-261.
- Gray, M.W., Butler, O. and Lang, B.F. (1997) *Science*, **283**, 1476-1481.
- Doolittle, W.F. (1999) *Science*, **284**, 2124-2128.
- Grosjean, H., Sprinzl, M. and Sierberg, S. (1995) *Biochimie*, **77**, 139-141.
- Bills, S.R., Hopper, A.K. and Martin, N.C. (1989) *Mol. Cell. Biol.*, **9**, 1611-1620.
- Sprinzl, M., Horst, C., Brown, M., Ioudovitch, A. and Steinberg, S. (1998) *Nucleic Acids Res.*, **26**, 148-153.
- Rose, A.M., Joyce, P.B., Hopper, A.K. and Martin, N.C. (1992) *Mol. Cell. Biol.*, **12**, 3652-3658.
- Dingwall, C. and Laskey, R.A. (1991) *Trends Biochem. Sci.*, **16**, 478-481.
- Ohno, M., Fomerod, M. and Martaj, I.W. (1998) *Cell*, **92**, 327-336.
- Rose, A.M., Belford, H.G., Shan, W.C., Greer, C.L., Hopper, A.K. and Martin, N.C. (1995) *Biochimie*, **77**, 45-53.
- Liu, J., Liu, J. and Stray, K.B. (1998) *Nucleic Acids Res.*, **26**, 3102-3108.
- Chiu, M.I., Mason, T.L. and Fink, G.R. (1992) *Genetics*, **132**, 987-1001.
- Akashi, K., Grandjean, O. and Small, I. (1998) *FEBS Lett.*, **431**, 39-44.
- Lund, E. and Dahlberg, J.E. (1998) *Science*, **282**, 2082-2085.
- Sarkar, S., Azad, A.K. and Hopper, A.K. (1999) *Proc. Natl Acad. Sci. USA*, **96**, 14366-14371.
- Schimmel, P. and Wang, C.-C. (1999) *Trends Biochem. Sci.*, **24**, 127-128.
- Mirande, M. and Waller, J.P. (1988) *J. Biol. Chem.*, **263**, 18443-18451.
- Mirande, M. (1991) *Prog. Nucleic Acid Res. Mol. Biol.*, **40**, 95-142.
- Kisilev, L.L. and Wolfson, A.D. (1994) *Prog. Nucleic Acid Res. Mol. Biol.*, **48**, 83-142.
- Yang, D.C.H. (1996) *Curr. Top. Cell. Regul.*, **34**, 101-135.
- Franklyn, C., Mosier-Forey, K. and Martin, S.A. (1997) *RNA*, **3**, 954-960.
- Yue, D., Maizels, N. and Weiner, A.M. (1996) *RNA*, **2**, 895-908.
- Wolfe, C.L., Lou, Y.C., Hopper, A.K. and Martin, N.C. (1994) *J. Biol. Chem.*, **269**, 13361-13366.
- Wolfe, C.L., Hopper, A.K. and Martin, N.C. (1996) *J. Biol. Chem.*, **271**, 4679-4686.
- von Heijne, G., Steppuhn, J. and Herrmann, R.G. (1989) *Eur. J. Biochem.*, **180**, 535-545.
- Cline, K. and Henry, R. (1996) *Annu. Rev. Cell. Dev. Biol.*, **12**, 1-26.
- Onofri, M., Hong, T., Nagelhus, T.A., Slupphaus, O., Lindino, T. and Krokan, H.E. (1998) *Nucleic Acids Res.*, **26**, 4611-4617.
- Muller-Weeks, S., Maniran, B. and Caradonna, S. (1998) *J. Biol. Chem.*, **273**, 21909-21917.
- Percival, K.J., Klein, M.B. and Burgers, P.M. (1989) *J. Biol. Chem.*, **264**, 2593-2598.
- Go, M. (1981) *Nature*, **291**, 90-93.
- de Souza, S.J., Lohg, M., Schenbach, L. and Gilbert, W. (1996) *Proc. Natl Acad. Sci. USA*, **93**, 14632-14636.